

Как достать всё что угодно со всего Интернета

Кучумов Илья,
руководитель разработки
Товарного Поиска



HighLoad++
2022

Яндекс

Обо мне



- + Руководитель отдела разработки Товарного Поиска в Яндексе
- + В прошлом руководил качеством поиска в Лавке и разработкой Турбо-страниц

Яндекс

айфон 13

Найти

Поиск

Картинки

Видео

Карты

Товары

Переводчик

Все

Мобильные телефоны

По популярности

1 Фильтры


Бренд X

Цена v

Линейка v

Операционная система v

Диагональ экрана v




4,8

Смартфон Apple iPhone 13 128 ГБ, тёмная ночь

49 990 ₽

Лучшая цена

ещё 33 предложения




4,8

Смартфон Apple iPhone 13 256 ГБ RU, сияющая звезда

64 990 ₽

Лучшая цена

ещё 5 предложений




4,8

Смартфон Apple iPhone 13 128 ГБ RU, сияющая звезда

54 978 ₽

Лучшая цена

ещё 5 предложений




4,8

Смартфон Apple iPhone 13 512 ГБ RU, сияющая звезда

94 980 ₽


Лучшая цена

ещё 5 предложений




4,8

Смартфон Apple iPhone 13 512 ГБ RU, синий




4,8

Смартфон Apple iPhone 13 512 ГБ RU, розовый



4,8

Смартфон Apple iPhone 13 128 ГБ, синий



4,8

Смартфон Apple iPhone 13 256 ГБ, тёмная ночь

Карты

Товары

Переводчик

Смартфон Apple iPhone 13 128 ГБ, тёмная ночь

4,8 948 отзывов 94% рекомендуют

Цвет товара: тёмная ночь

Версия: Ростест (ЕАС) для других стран

Конфигурация памяти: 128 ГБ 256 ГБ 512 ГБ

Коротко о товаре

экран 6.1" (2532x1170) OLED 60 Гц

встроенная память 128 ГБ

[Все характеристики](#)

Предложение спонсора Реклама

re-store.ru 75 990 ₽ В магазин

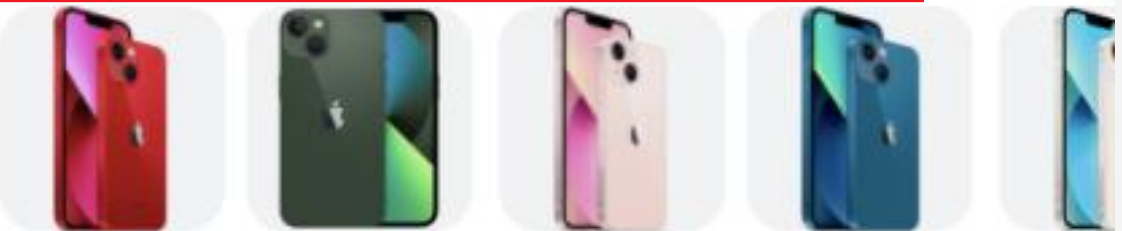
Цены в магазинах

ant-shop.ru 4,6 732 отзыва 49 990 ₽ Лучшая цена В магазин

berudevice.ru 4,5 Нет отзывов 50 500 ₽ В магазин

4,8 948 отзывов 94% рекомендуют

Цвет товара: тёмная ночь



Конфигурация памяти:

128 ГБ 256 ГБ 512 ГБ

Цены в магазинах

ant-shop.ru 4,6 732 отзыва

Про что доклад

- + Про построение базы цен Товарного Поиска
- + Мы проверили, все работает и на других срезах
- + Ключевая функциональность, высокие требования правильности
- + Цены быстро меняются, невозможно кэшировать

План доклада

- 01 Общая архитектура базы Товарного Поиска
- 02 Подробно про парсинг и скачивание
- 03 Другие удачные применения парсинга

Что такое построение базы цен

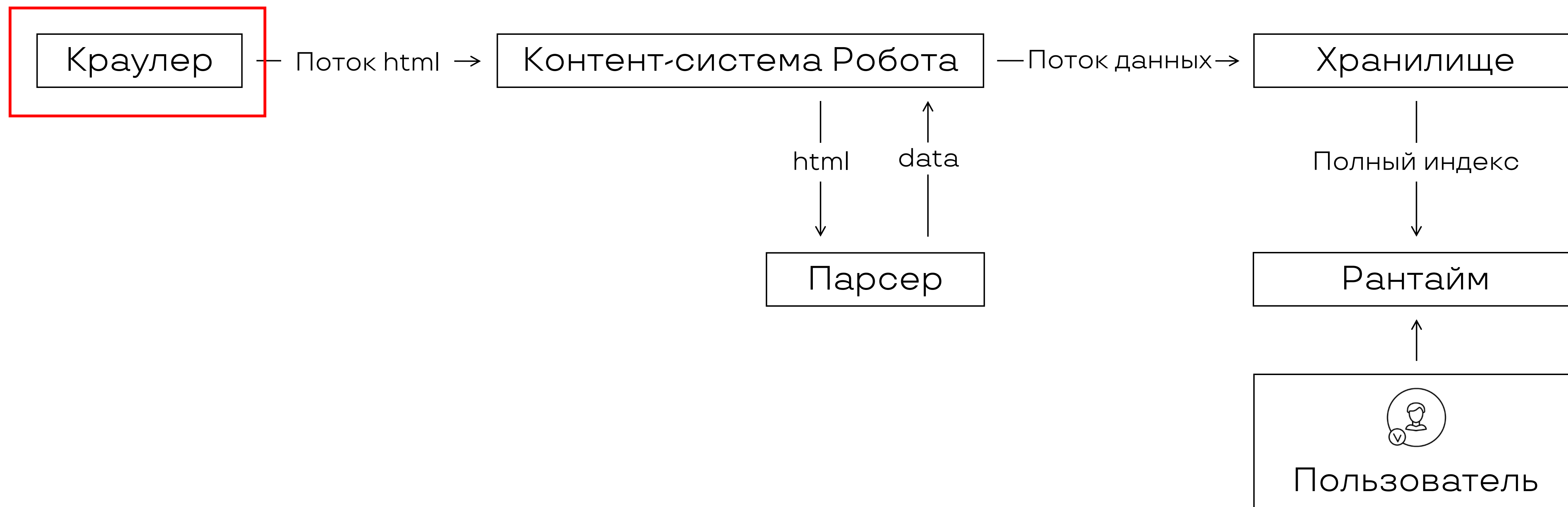
20 млрд
товаров

1,5 млн
интернет-магазинов

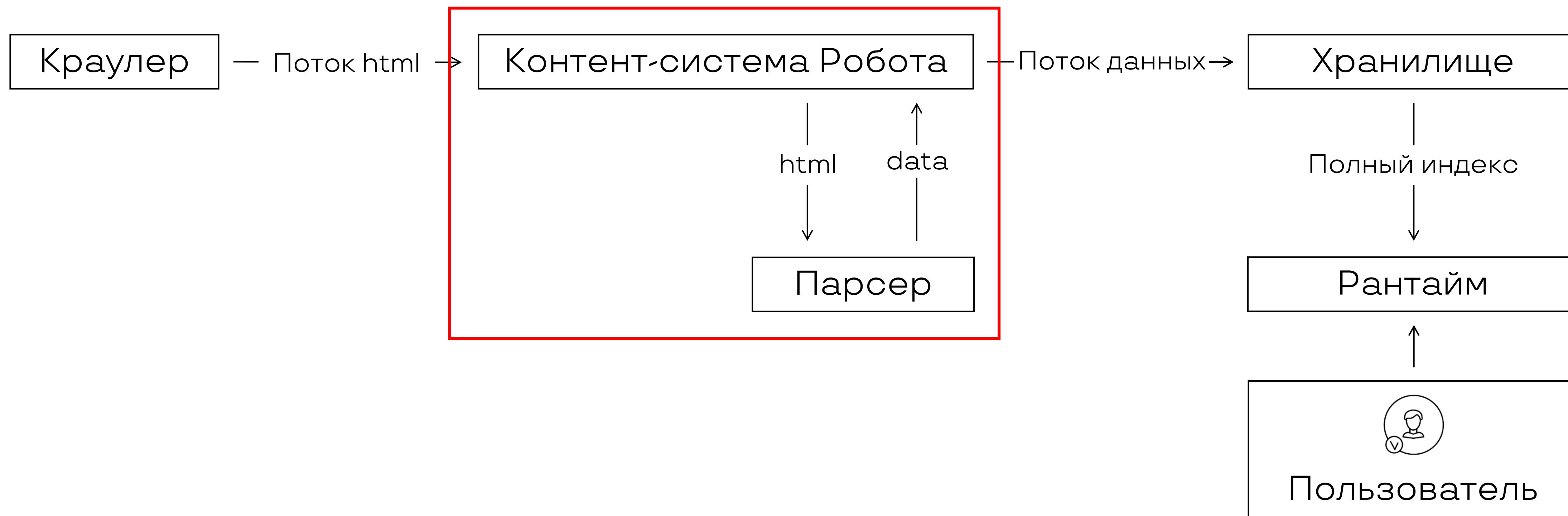
~11% случайных товаров меняют
цену каждые **24** часа

Ошибка цены — это плохо

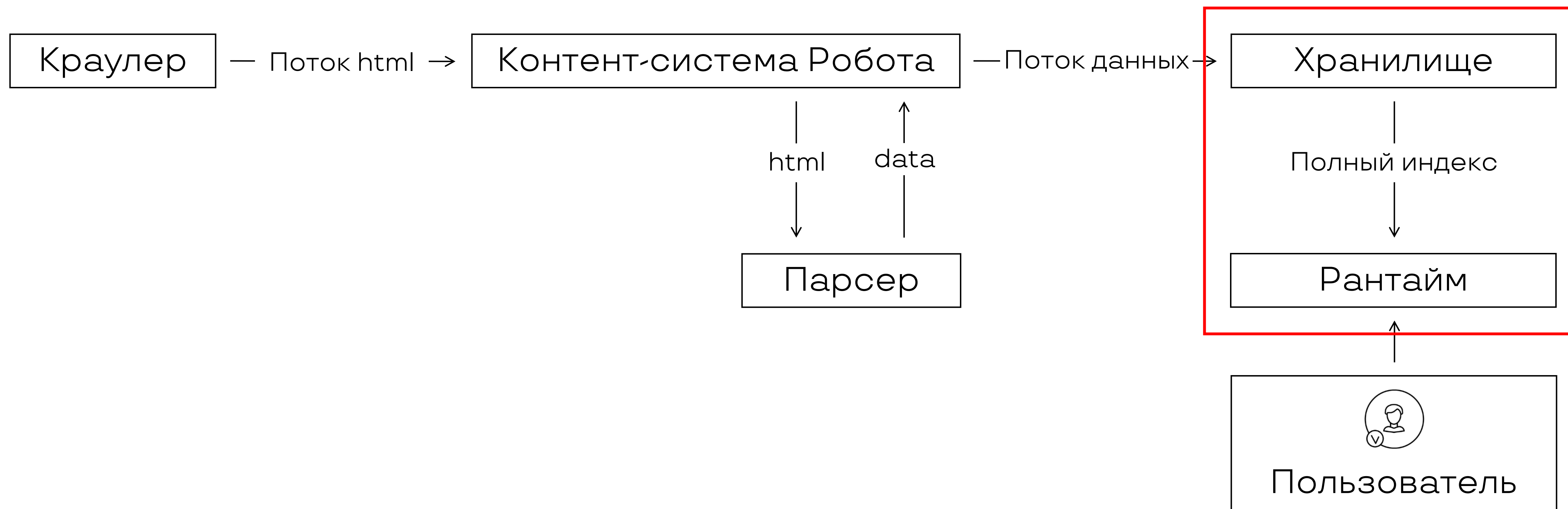
Общая архитектура



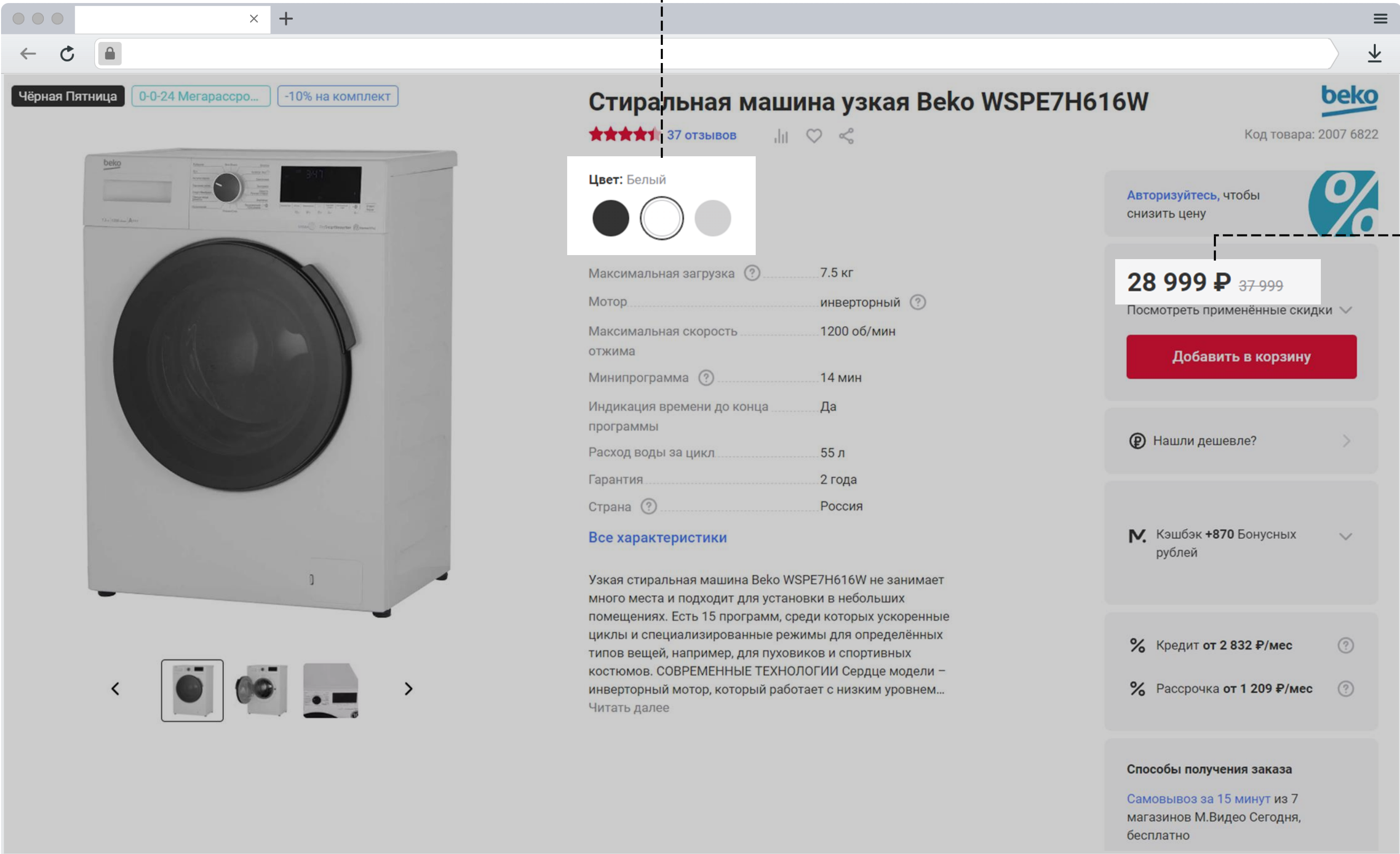
Общая архитектура



Общая архитектура



Что такое парсинг



Цвет: Белый



28 999 ₪ 37 999

Что такое парсинг

Электросамокат Ninebot KickScooter D18U, до 100 кг, черный

4.6 124 отзыва Характеристики 41 вопрос дизайн, качество сборки 23 370 человек интересовались за 2 месяца

Следить за снижением цены

В избранное

Сравнить

Эксклюзив

ЕЩЕ 17

Коротко о товаре

Максимальная нагрузка

100 кг

Запас хода

до 18 км без подзарядки

Максимальная скорость

25 км/ч

Мощность двигателя

250 Вт

Вес

14.9 кг

Особенности конструкции

амортизатор, складной

Материал рамы

алюминиевый сплав

Материал колес

резина

Аккумулятор

5100 мА·ч

Тип тормоза

ручной, комбинированный

Подробнее

Задать вопрос о товаре

Все товары Ninebot

СУПЕР ЦЕНЫ 19.21 октября

19 990 ₽

32 990 ₽

8 248 ₽ x 4 сплитом

660 баллов на Плюс

Самовывоз завтра, 19 октября — бесплатно

Курьером завтра, 19 октября — бесплатно

Доставка Яндекс со своего склада

Оплата онлайн

1 товар в корзине

Маркет

2 мес

4 мес

6 мес

8 248 ₽ сегодня

и 24 742 ₽ потом

Оформить

18 окт

1 ноя

15 ноя

29 ноя

8 248 ₽

8 248 ₽

8 248 ₽

8 246 ₽

СУПЕР ЦЕНЫ 19.21 октября

19 990 ₽

32 990 ₽

12

Яндекс

HL

HighLoad++

2022

Что такое парсинг

HTML →

```
{  
  "title": "Стиральная машина",  
  "price": 32999,  
  "available": true,  
  "main_image": "https:cdn.shop.ru/wash.jpg",  
  "brand": "Самый лучший"  
}
```

А можно ли просто?

Микроразметка

```
<meta itemprop="availability" content="https://schema.org/InStock" />  
<meta itemprop="priceCurrency" content="RUB" />  
<meta itemprop="price" content="82990" />
```

- + Быстро реализовать
- + Веб-мастера напрямую влияют
- + Покрытие — 50% товаров
- + Невозможно влиять со стороны парсера

Парсеры на селекторах

```
<div class="ProductCartFixedBlock__price"> ←  
  <span class="_current-price js--_current-price ">  
    5990  
  </span>  
  <span class="ProductPrice__rouble __rouble">  
    &#x20bd;  
  </span>  
</div>
```

Парсеры на селекторах

Новый парсер

Предпросмотр

Селекторы

Нажмите селекторы с помощью Watson. Вы можете выделять блоки, коллекции и указывать дополнительную обработку контента.

Блок

price

.price-block__final-price

Вывод

Извлечь содержимое

Sanitize / Исключить теги

Микроразметка

Парсинг микроразметки страницы.

Спарсить микроразметку

Не парсить микроразметку

Поддержка Watson

Скрыть расширение

Apple / Смартфон Apple iPhone 11 128GB

РАССРОЧКА ОТ 0-0-6

☆☆☆☆☆ 212 отзыва

Артикул: 86123929

Купили более 1 400 раз

Цвет: черный

Питание

Время работы в режиме разговора (max)

Емкость аккумулятора

Мультимедийные возможности

Фронтальная камера (млн. пикс.)

12 Мп

Количество мп основной камеры

12 Мп

История цены

от 43 800 Р до 54 99

Селекторы

div.product-page__price-block:nth-child(1) > div:nth-child(8) > div:nth-child(3) > p:nth-child(1) > span:nth-child(1) > ins:nth-child(3)

.price-block__final-price

ins

Блок

price

.price-block__final-price

Вывод

Извлечь содержимое

Sanitize / Исключить теги

Селекторы

div.product-page__price-block:nth-child(1) > div:nth-child(8) > div:nth-child(3) > p:nth-child(1) > span:nth-child(1) > ins:nth-child(3)

.price-block__final-price

ins

price

.price-block__final-price

Вывод

Извлечь содержимое

Sanitize / Исключить теги

.price-block__final-price

50 040 Р

54 990 Р

В кредит от 5 004 Р

Добавить в корзину

20-22 октября доставит Wildberries со склада Коледино WB

Попытка накликать парсеры

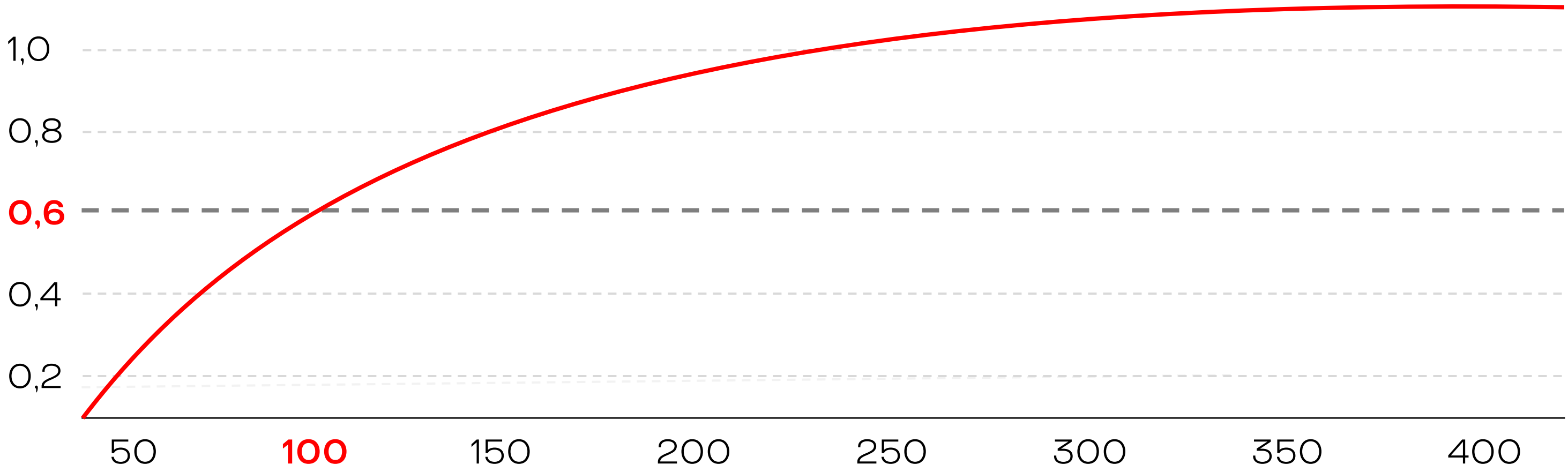
Apple iPhone 13 Pro Max 256GB Dual SIM Silver купить...

best-magazin.com > Apple iPhone 13 Pro Max 256GB Dual SIM Silver (... ★ 4,9

Купить Айфон 13 Про Макс 256 ГБ Dual SIM в Бест-магазин. ... Процессор iPhone 13 Pro Max Dual SIM (256GB, Silver, 2 nano SIM) имеет 2 высокопроизводительных ядра (тактовая частота до 3,23 ГГц) плюс 4 ядра эффективности. Читать ещё

99 990 ₽





Попытка накликать парсеры

10 разработчиков

5 пицц

- + Кола и чай
- + Большая переговорка

Мы написали ~~100~~ **80** парсеров за встречу

~~80~~ **70** отправили в эксперимент
(через месяц осталось 60)

Парсеры на селекторах

- + Возможность влиять на качество со стороны парсера
- + Быстро распарсить конкретный сайт
- + 1.5 млн. парсеров = около 130 рабочих человеко-лет накликивания
- + Нужна постоянная поддержка разработчиков

Парсеры краудом

01



Берем
конкретный
сайт

02



Толокер
создает
селекторы

03



Другие толокеры
проверяют
качество
на семпле
страниц

04

Если всё
хорошо —
добавляем в
стейт
парсера

Почему неидеальное решение?

+ Парсеры независимые

1.5 млн. парсеров = 1.1 млн. долларов 0,75 долларов/штука

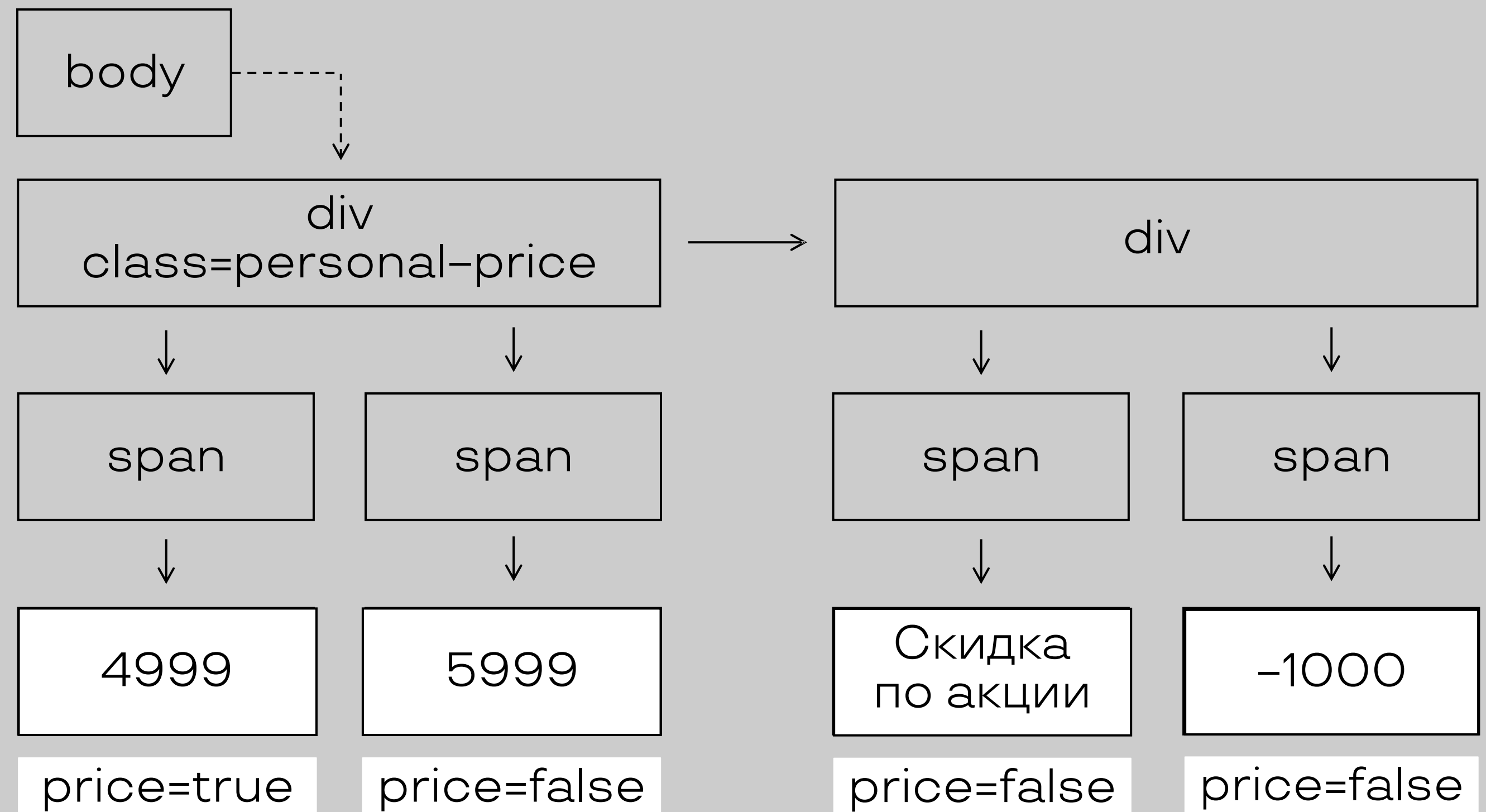
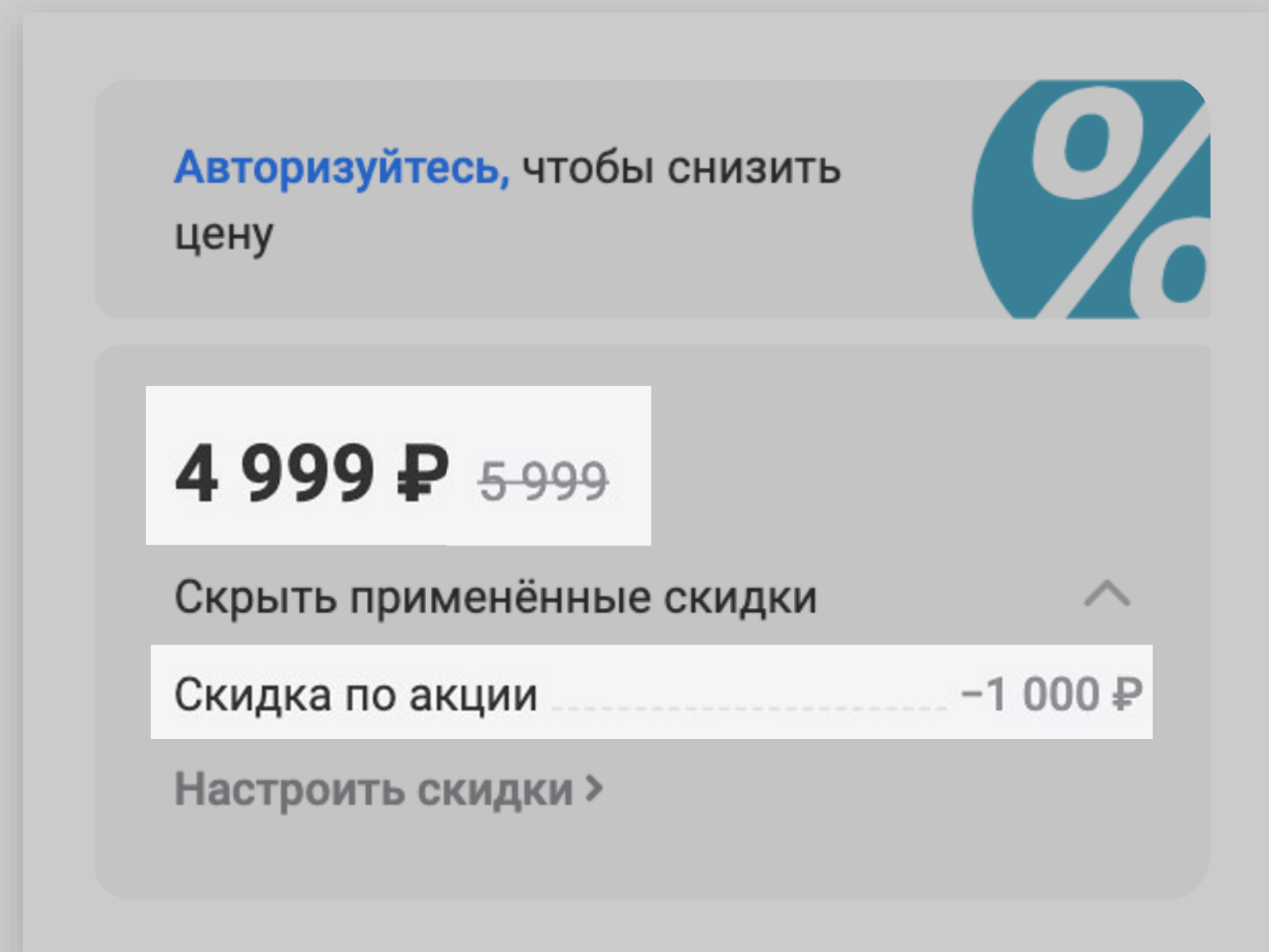
30% случайных парсеров
сломаются за месяц

- ! Стоимость проверки 750 тыс. долларов (еженедельно)
- ! Расходы не делятся на другие срезы

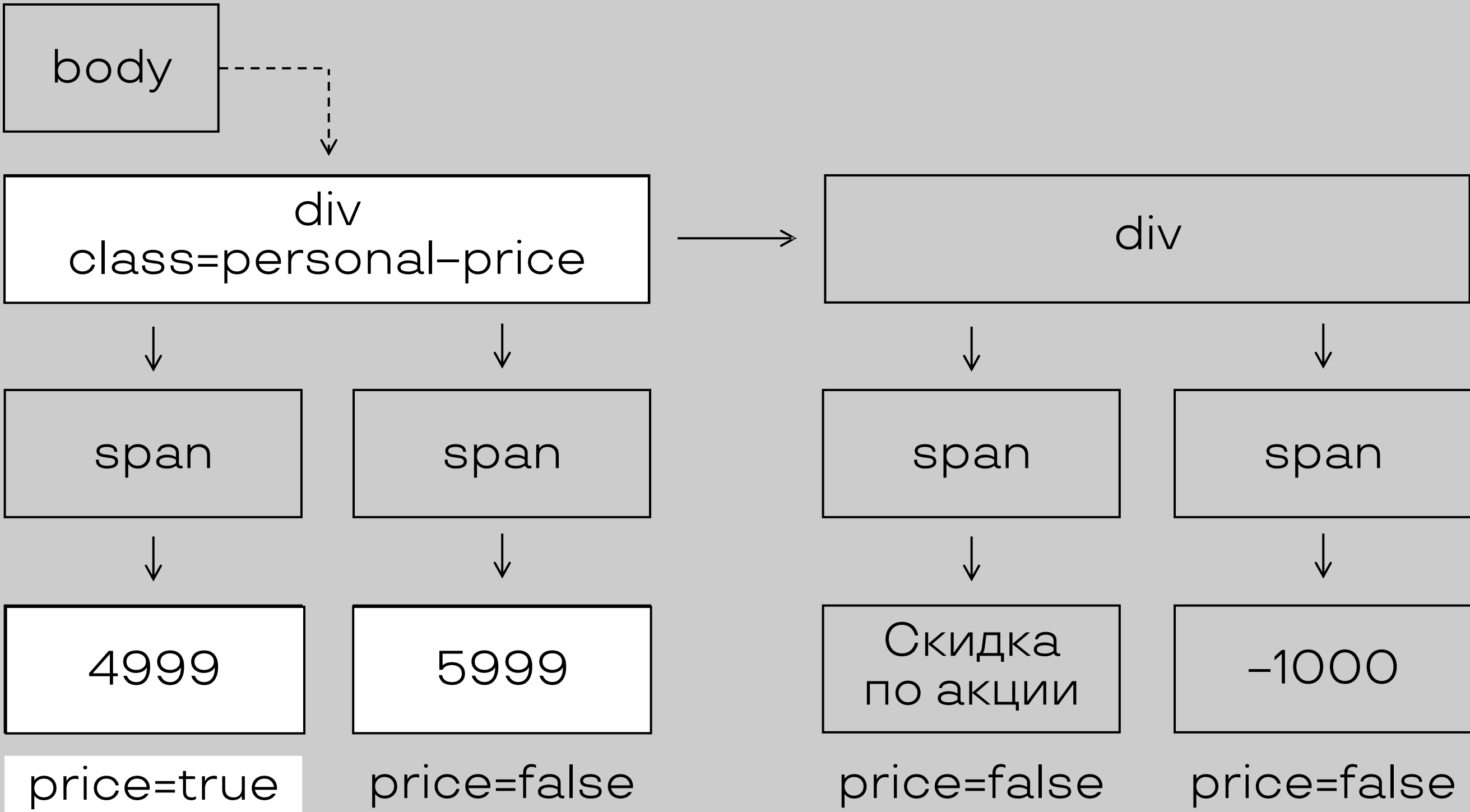
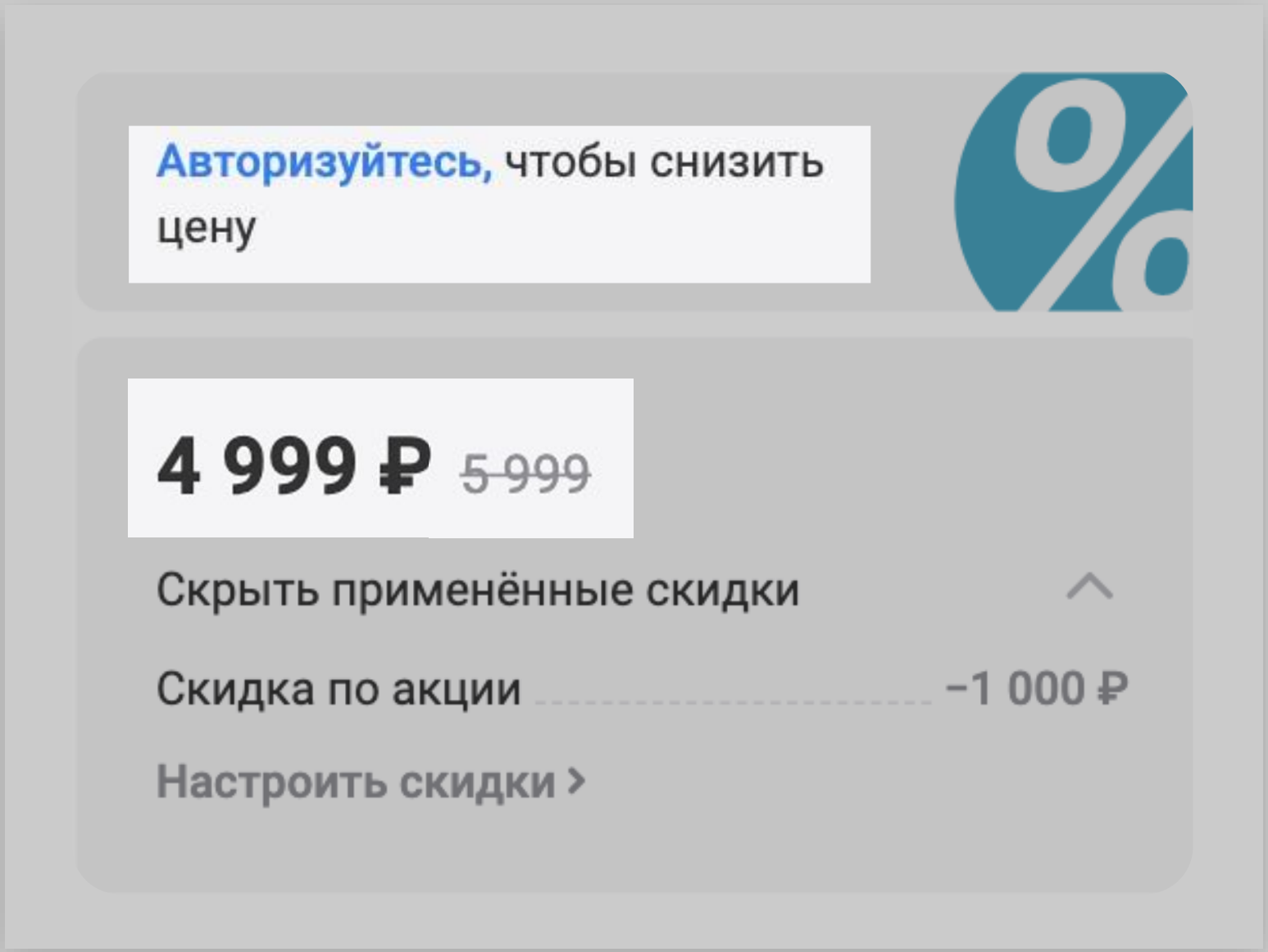
Саммари простых методов

- + Простых рабочих эвристик по контенту нет
- + Микроразметка — среднее покрытие, но быстро запустить
- + Селекторы — хорошее покрытие, но дорого поддерживать

Классификация вершин DOM-дерева



Текстовый фактор

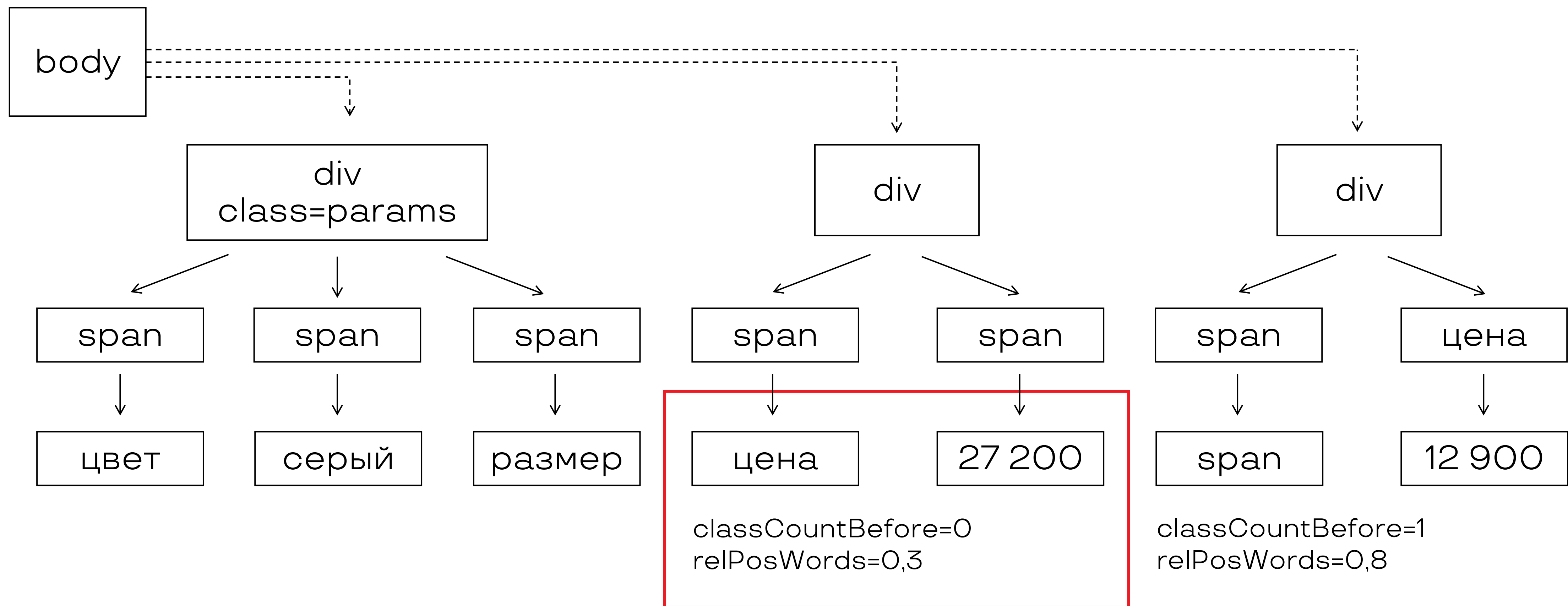


Text = Авторизуйтесь, чтобы снизить цену | 4 999 || class=personal-price | 5999

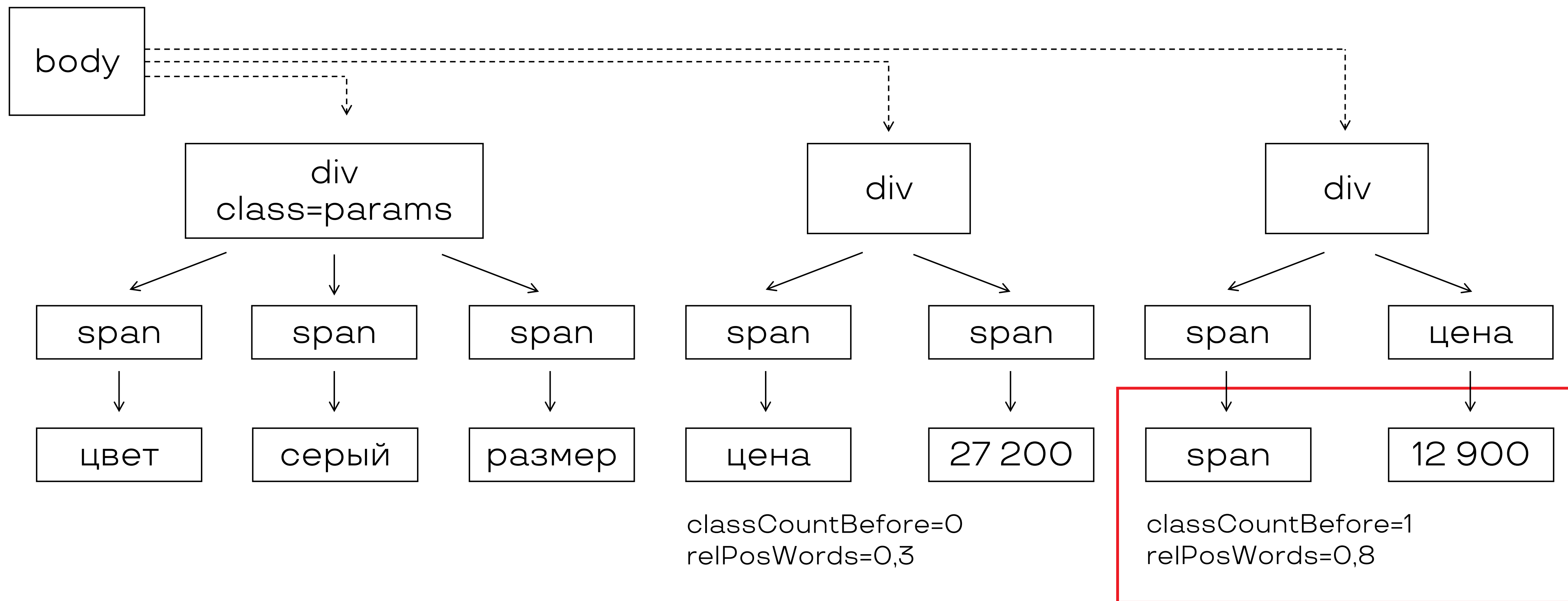
Контекст вершины в дереве

- + Тег родителя, дедушки, ...
- + Сколько параграфов (ссылок, картинок, ...) у родителя, дедушки?
- + А сколько еще таких же вершин в дереве?
- + Количество вершин у родителя, дедушки, ...

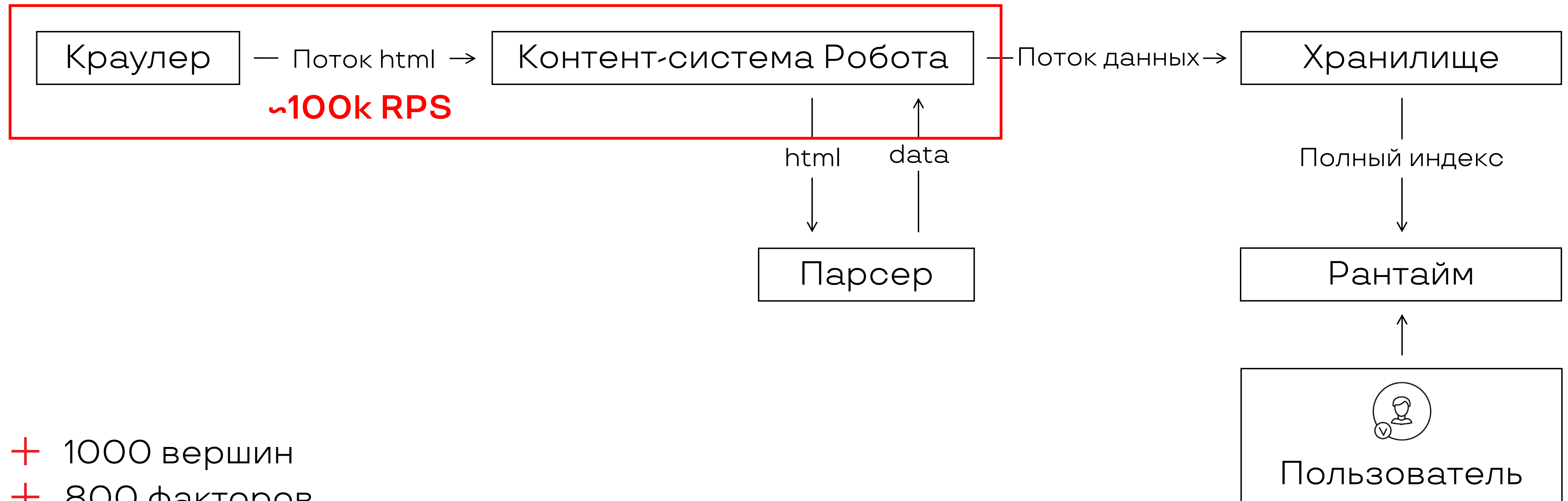
Позиция в дереве



Позиция в дереве

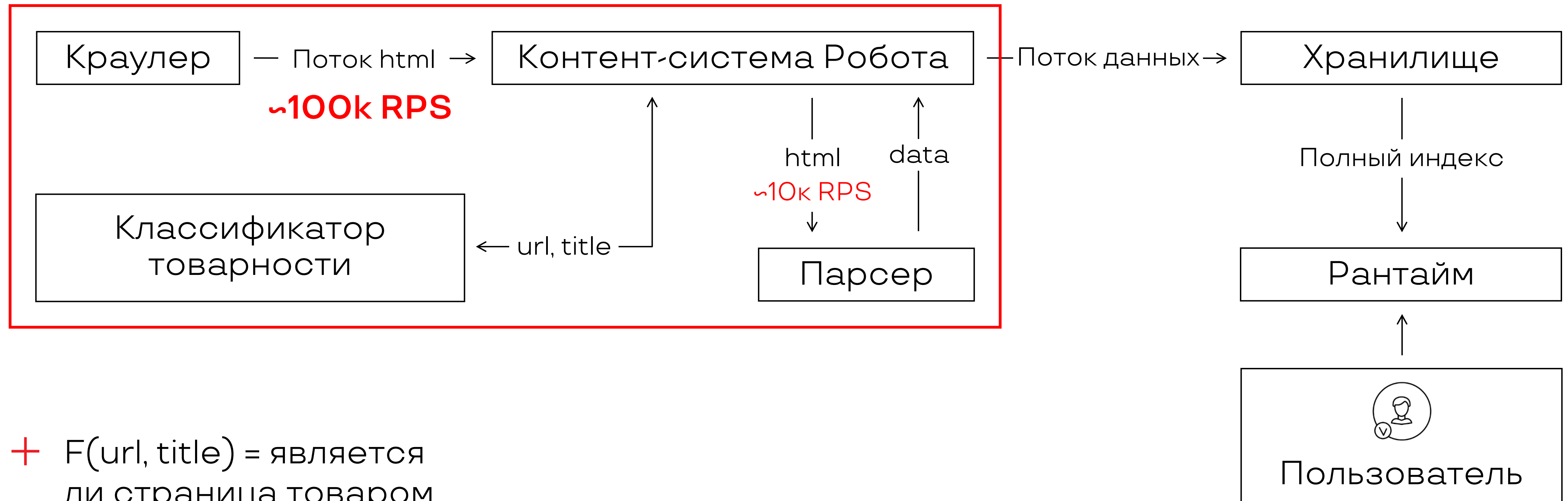


Необходимость оптимизаций



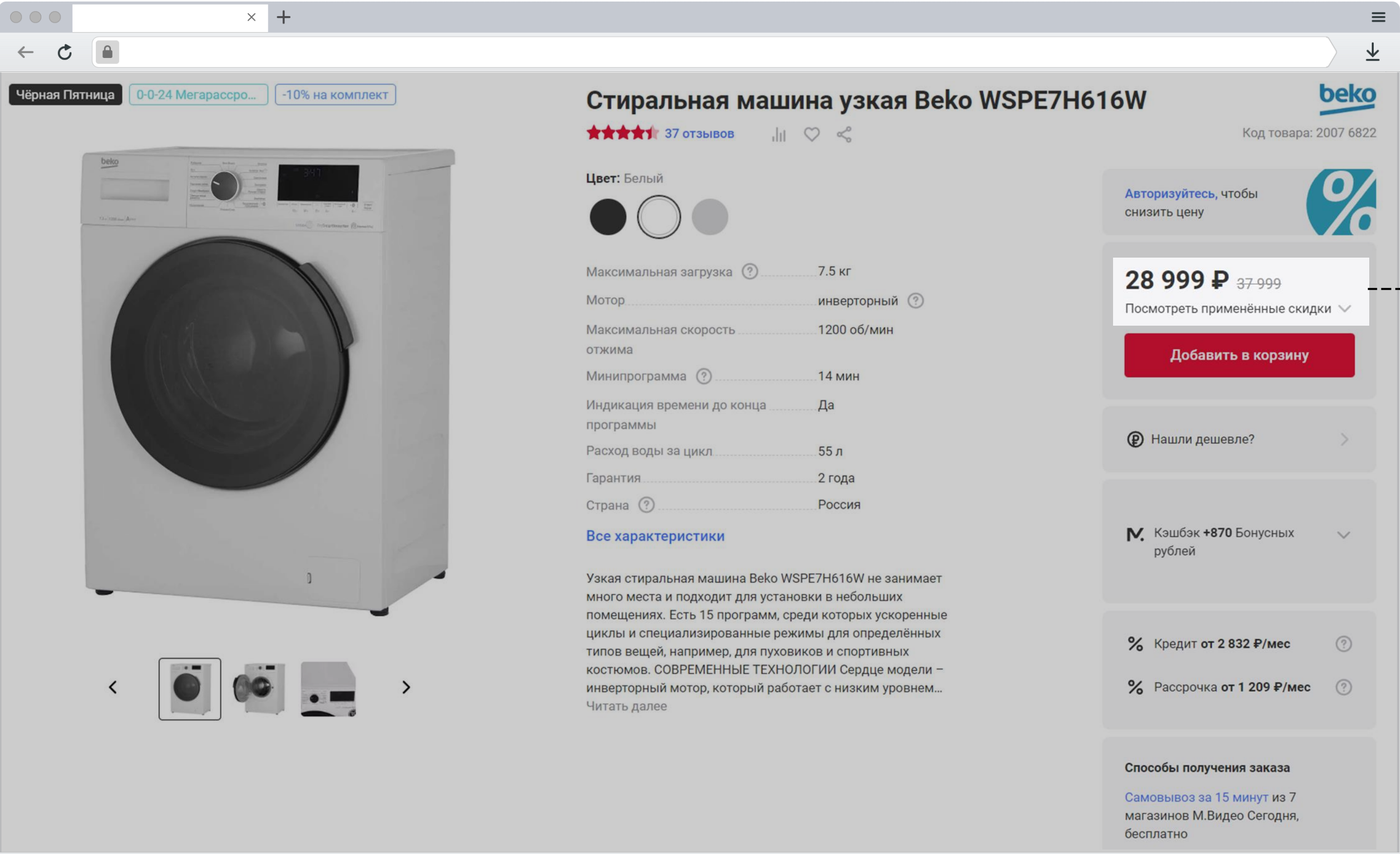
- + 1000 вершин
- + 800 факторов
- + ~ 10к CPU

Сужение потока до товарных страниц

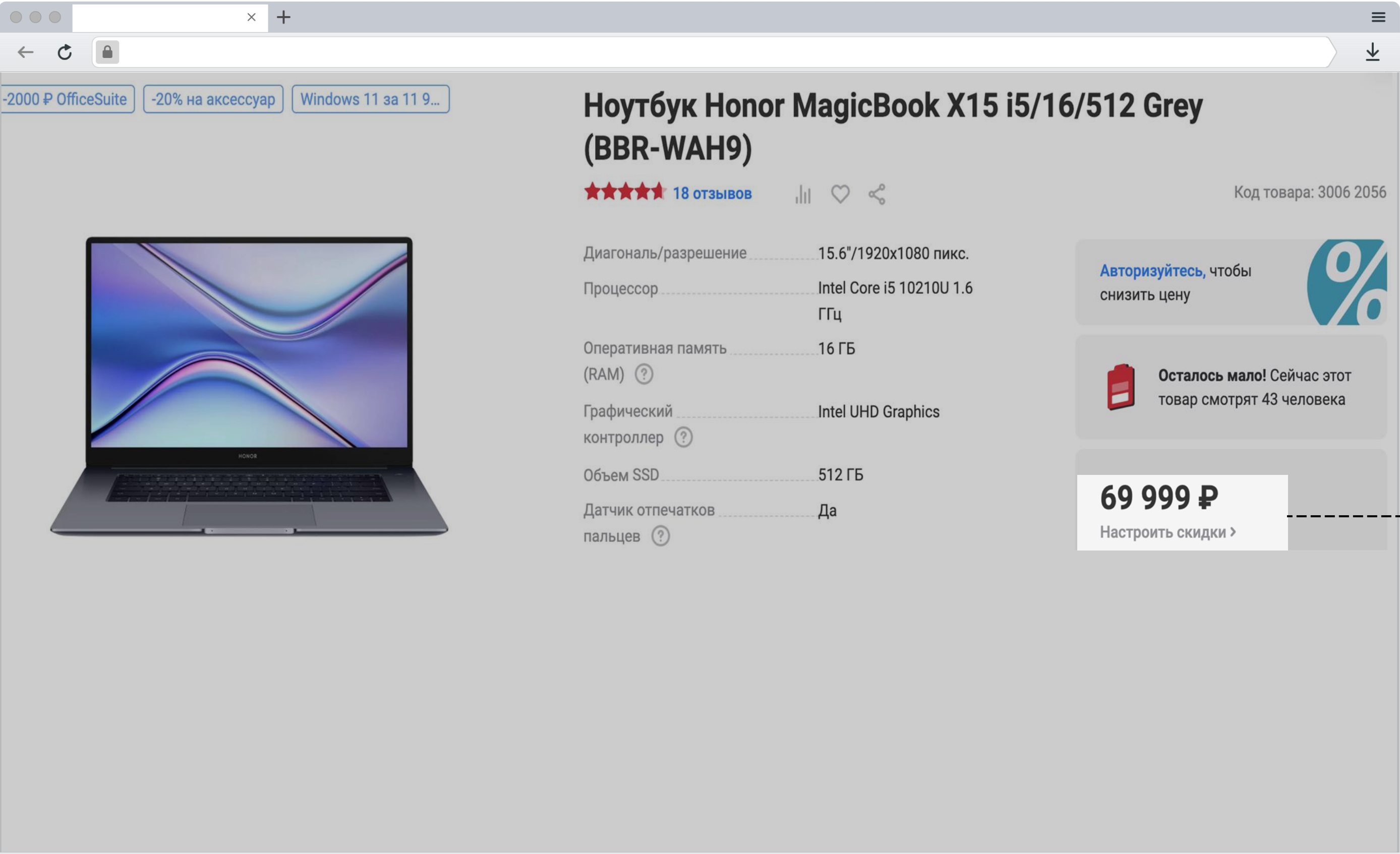


- + $F(url, title)$ = является ли страница товаром
- + Легкая DSSM'ка

Похостовое кэширование



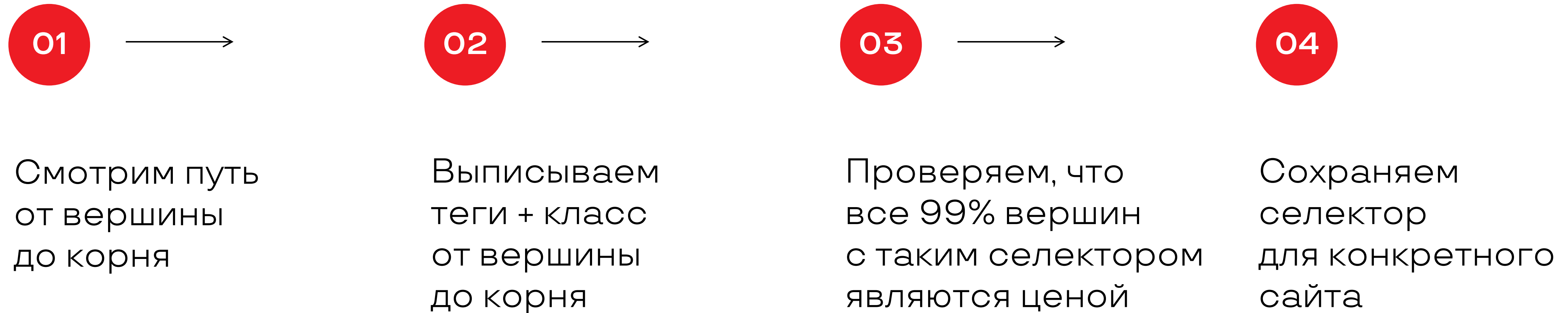
Похостовое кэширование



69 999 ₽
Настроить скидки >

Давайте генерировать селекторы моделью

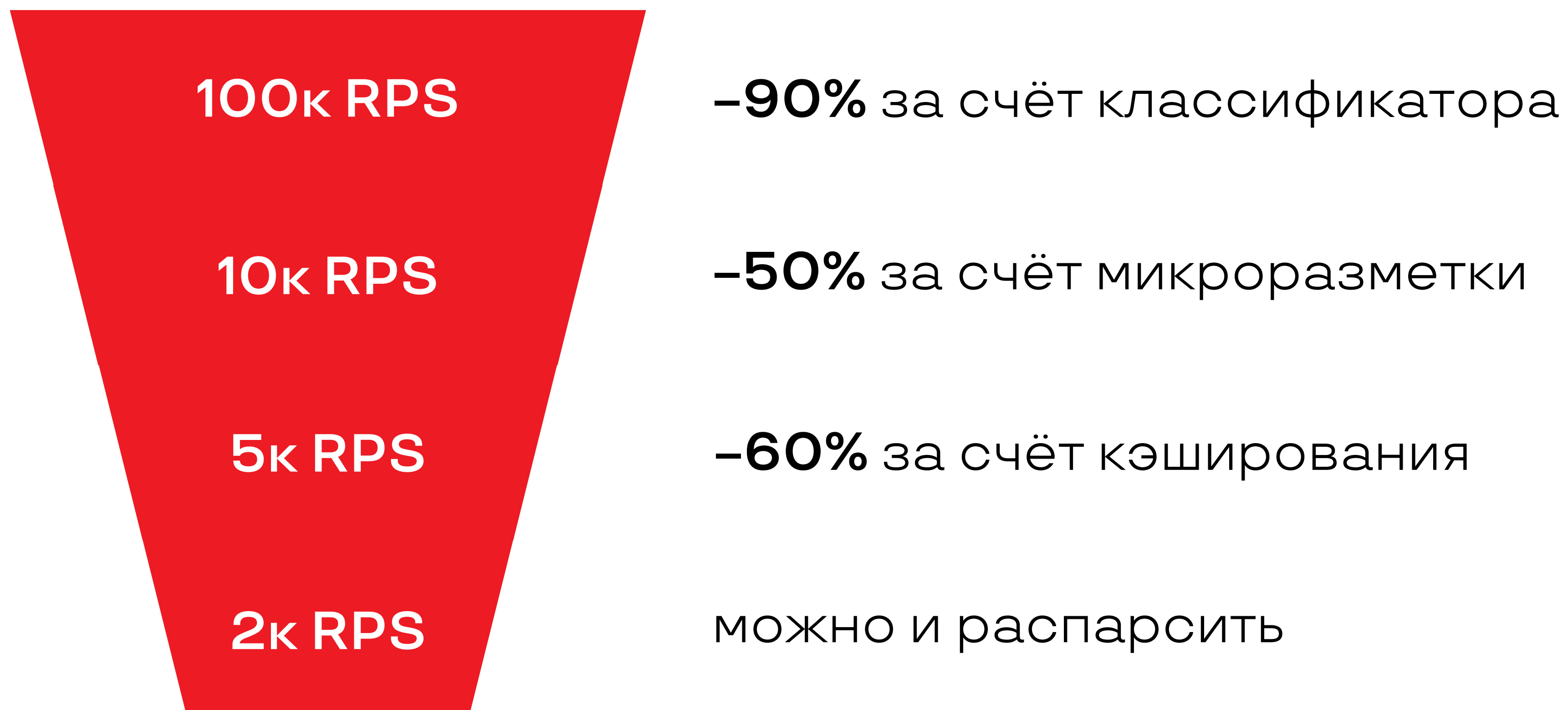
Алгоритм построения селекторов



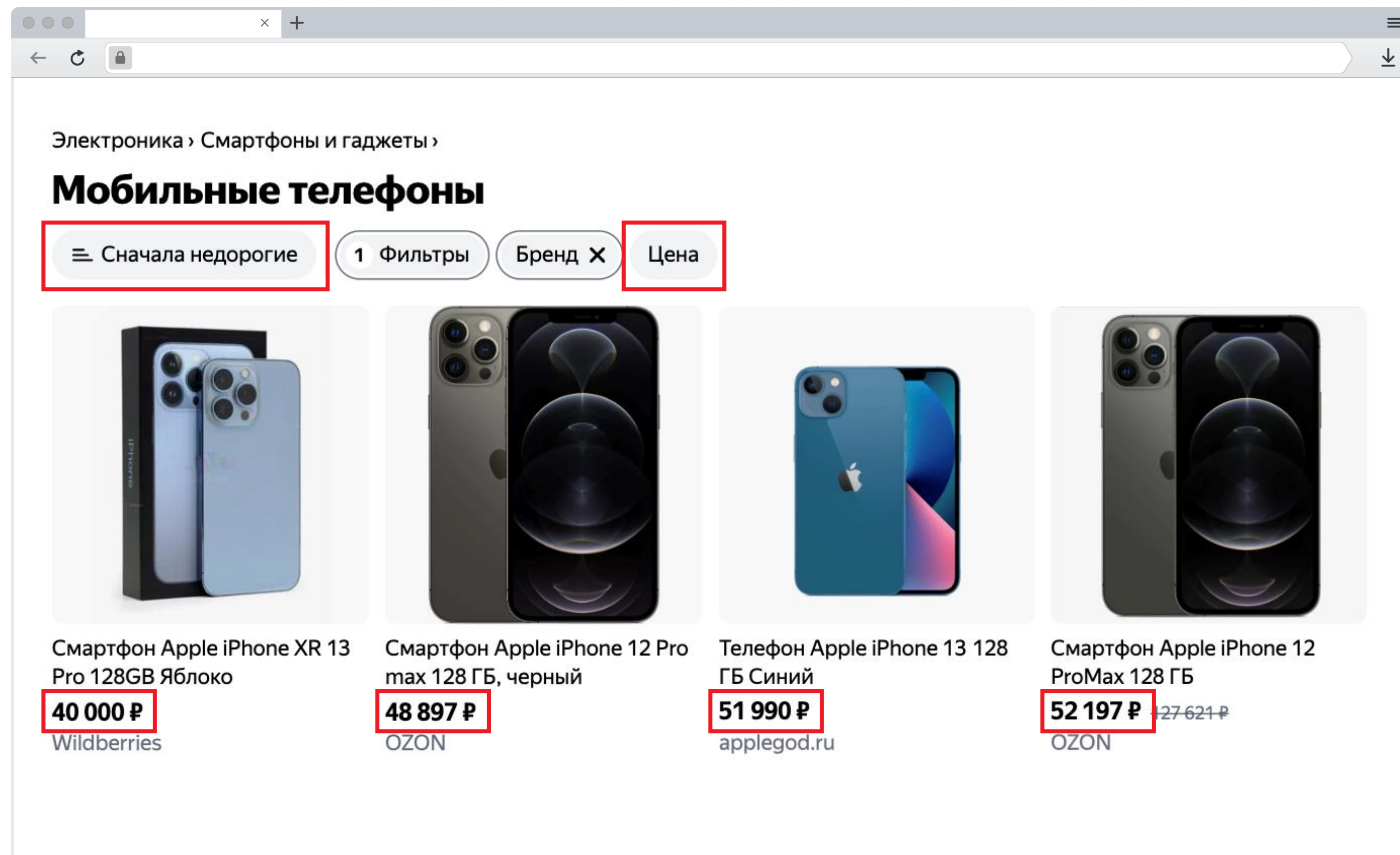
Как используем

- + Для каждого хоста (host, css_selector) → answer
- + Файлик с миллионом записей
- + Регулярно перестраиваем без участия крауда
- + Экономия 60% железа

Как режется поток



Почему актуальность цен важна



**Давайте будем обходить базу
раз в день (неделю?)**

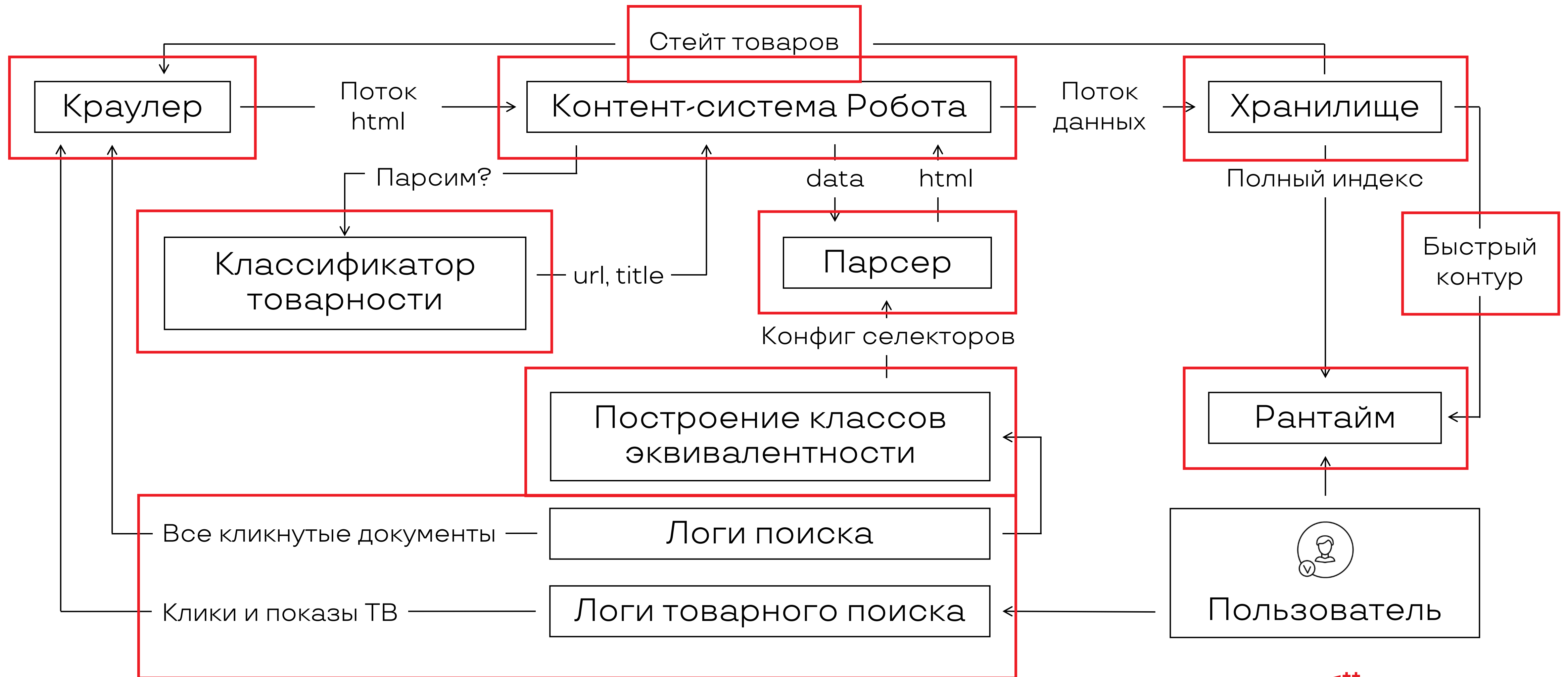
Поддержка актуальности

- + Round robin не работает, прокачка всей базы — больше месяца
- + Похостовые лимиты — серьёзная проблема
- + 10М за день — это 116 RPS на конкретный сайт

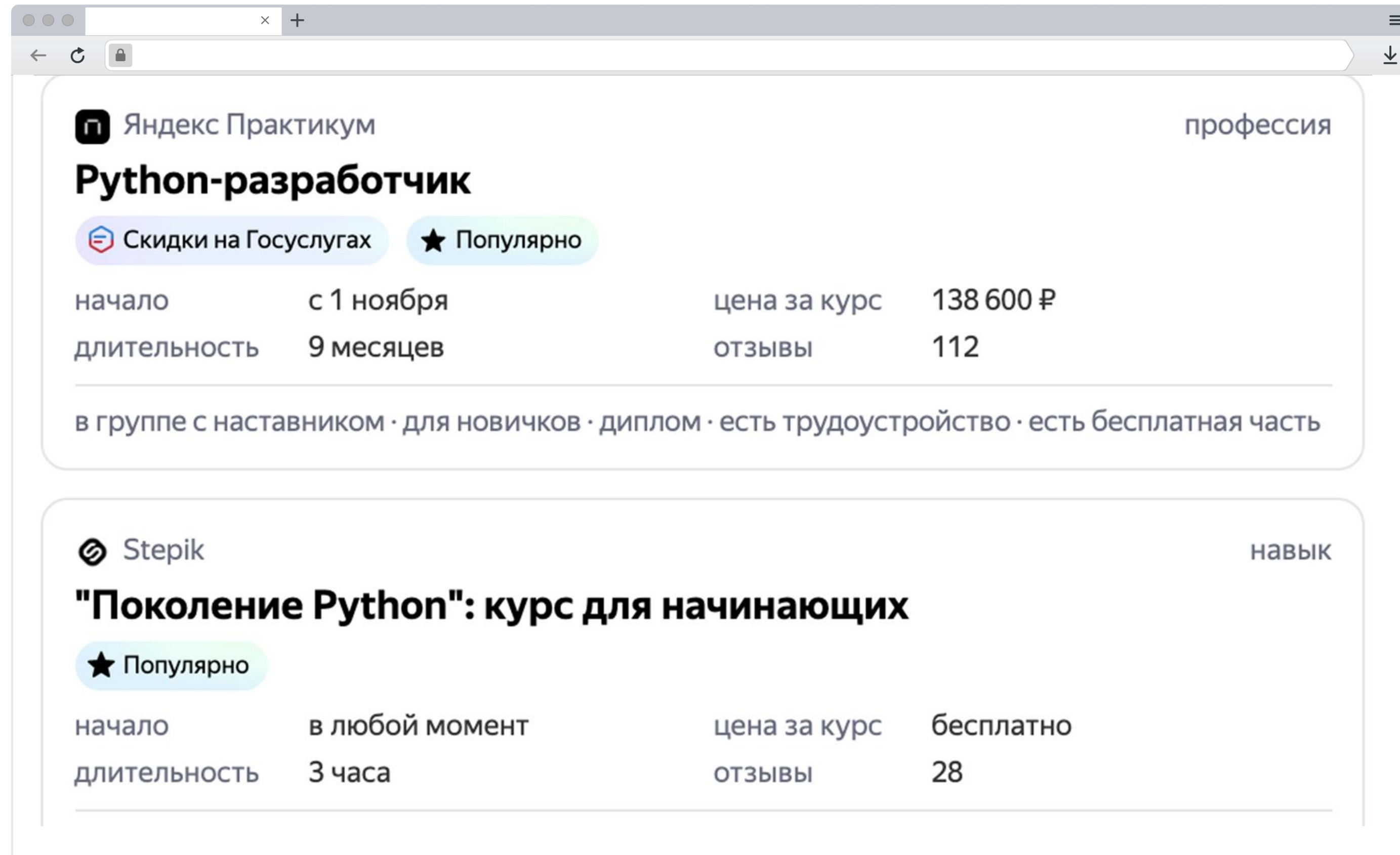
Эвристики скачивания

- + Показы в поиске обходим всегда — покрываем популярные
- + По (title, category, offer_age) предсказываем вероятность изменения цены с последнего скачивания
- + Качаем самые вероятно изменившиеся
- + Подбираем фильтр по возрасту, если нужно

Финальная схема



Что можно парсить ещё



The screenshot shows a web browser window with two course listings. The first listing is from Yandex Practicum for a 'Python-developer' profession. It includes details like start date (November 1st), duration (9 months), price (138,600 RUB), and reviews (112). The second listing is from Stepik for a 'Python Generation' course for beginners, which is free and has a duration of 3 hours. Both listings are marked as 'Popular'.

Platform	Course Name	Category	Start Date	Duration	Price	Reviews	Additional Info
Яндекс Практикум	Python-разработчик	профессия	с 1 ноября	9 месяцев	138 600 ₽	112	Скидки на Госуслугах, Популярно, в группе с наставником, для новичков, диплом, есть трудоустройство, есть бесплатная часть
Stepik	"Поколение Python": курс для начинающих	навык	в любой момент	3 часа	бесплатно	28	Популярно

Сложные характеристики товаров

Куртка мужская МОД 12032, темно-серый, 48/182

Цвет товара:

темно-серый

черный

Размер:

(48)182-96-80

(50)182-100-84

Коротко о товаре

цвет товара

темно-серый

сезон

зима

состав

полиамид 100%

силуэт

прямой

застежка

молния, кнопки

детали

карманы, капюшон

Резюме

- + Несколько подходов к парсингу структурированных данных
- + Как поддерживать актуальность базы
- + Архитектура базы Товарного Поиска
- + Различные сферы применения алгоритмов парсинга

Резюме

- † Несколько подходов к парсингу структурированных данных
- † Как поддерживать актуальность базы
- † Архитектура базы Товарного Поиска
- † Различные сферы применения алгоритмов парсинга

Кучумов Илья

руководитель разработки Товарного Поиска

@IlyaYndx

Обратная связь и комментарии
по докладу по ссылке



HighLoad⁺⁺
2022

Яндекс